MEAN SQUARE ERROR BEHAVIOR FOR PREDICTION

IN LINEAR REGRESSION MODELS

BY

ALAN E. GELFAND

TECHNICAL REPORT NO. 427

MARCH 7, 1990

PREPARED UNDER CONTRACT

N00014-89-J-1627    (NR-042-267)

FOR THE OFFICE OF NAVAL RESEARCH

DEPARTMENT OF STATISTICS

STANFORD   UNIVERSITY

STANFORD, CALIFORNIA

DTIC
ELECTE
APR 03 1990
E

90  04  02   121

MEAN SQUARE ERROR BEHAVIOR FOR PREDICTION

IN LINEAR REGRESSION MODELS

BY

ALAN E. GELFAND

TECHNICAL REPORT NO. 427

MARCH 7, 1990

DEPARTMENT OF STATISTICS

STANFORD   UNIVERSITY

STANFORD,  CALIFORNIA

DTIC
COPY
INSPECTED
4

# MEAN SQUARE ERROR BEHAVIOR FOR PREDICTION IN LINEAR REGRESSION MODELS

Alan E. Gelfand

## ABSTRACT

For the problem of individual prediction in linear regression models, that is, estimation of a linear combination of regression coefficients, mean square error behavior of a general class of adaptive predictors is examined.

## 1. INTRODUCTION

Suppose the usual linear regression model with fixed regressors, $Y = X\beta + \epsilon$, $Y_{n \times 1}$, $X_{n \times p}$ full rank, $\beta_{p \times 1}$ and $\epsilon_{n \times 1} \sim (0, \sigma^2 I)$. Let $\hat{\beta}_{LS} = (X^T X)^{-1} X^T Y$ denote the ordinary least squares estimator of $\beta$. At a new vector of predictor values, $X_0$, we seek to estimate $X_0^T \beta$. Using mean square error as a criterion, results of Cohen (1965) show that if $\epsilon$ is normally distributed, $\alpha X_0^T \hat{\beta}_{LS}$ is an admissible estimator of $X_0^T \beta$ for $0 \leq \alpha \leq 1$, e.g., the UMVU predictor is admissible. In fact, a predictor of the form $\ell^T Y$ is admissible for $X_0^T \beta$ iff $(2\ell - X(X^T X)^{-1} X_0)^T (2\ell - X(X^T X)^{-1} X_0 \leq X_0^T (X^T X)^{-1} X_0$.

In the sequel, we study the MSE under normality of predictors of the form $X_0^T \hat{\beta}_C$ where

$$\hat{\beta}_C = C\hat{\beta}_{LS} + (I - C)\beta^\star \qquad (1)$$

1

C a matrix usually data dependent and $\beta*$ a specified vector. Such $\hat{\beta}_C$ include most alternatives to $\hat{\beta}_{LS}$ discussed in the literature. Earlier work in this direction appears in Baranchik (1964) and Radhakrishnan (1970).

## 2. NOTATION AND MOTIVATION

To simplify matters, we convert to canonical form. Let $\hat{\alpha} = P\hat{\beta}_{LS}$, P orthogonal such that $P(X^T X)^{-1} P^T = D^{-1}$, D diagonal with diagonal elements $d_i$. Define $\alpha = P\beta$, $\ell = P\underline{X}_0$ and for convenience set $\beta* = 0$. For the moment assume $\sigma^2$ known. Our problem now is to estimate $\theta = \ell^T \alpha$ given $\hat{\alpha} \sim N(\alpha, \sigma^2 D^{-1})$ wishing to do well near $\theta = 0$. Let $U = \ell^T \hat{\alpha}$, $Z = \hat{\alpha}^T D\hat{\alpha}$, $q = \ell^T D^{-1}\ell$, $V = Z - U^2/q$, $\lambda = \alpha^T D\alpha$ and $\zeta = \lambda - \theta^2/q$. Then, U, V are independent, $U \sim N(\theta, \sigma^2 q)$, $V \sim \sigma^2 \chi^2_{p-1}(\zeta/\sigma^2)$.

Consider a general adaptive predictor $\delta(\hat{\alpha})$ of the form

$$\delta(\hat{\alpha}) = \Sigma h_i(\hat{\alpha}) \ell_i \hat{\alpha}_i. \tag{2}$$

Most predictors of $\theta$ discussed in the literature are special cases of (2). Apart from the LS predictor, U, we have:

i) A class of predictors given in Thompson (1968)

$$T_m = \frac{U^2}{U^2 + m\sigma^2 q} U, \text{ m a known constant, i.e., } h_i(\hat{\alpha}) = \frac{(\hat{\alpha}^T \ell)^2}{(\hat{\alpha}^T \ell)^2 + m\sigma^2 q}.$$

ii) A class of predictors given in Mehta and Srivastava (1971)

$$MS_{b_1, b_2} = (1 - b_1 e^{-b_2 U^2/\sigma^2 q})U, \; 0 < b_1 < 1, \; b_2 > 0, \; b_1, b_2 \text{ known},$$

$$\text{i.e., } h_i(\hat{\alpha}) = 1 - b_1 \exp(-b_2(\hat{\alpha}^T \ell)^2/\sigma^2 q).$$

iii) A predictor arising from the James-Stein estimator adapted for unequal variances (Sclove 1968)

$$JS_c = (1 - \frac{c\sigma^2}{Z})U, \text{ c known usually taken equal to p - 2.}$$

A positive part adjustment should be applied so that $h_i(\hat{\alpha}) = [1 - c\sigma^2(\hat{\alpha}^T D\hat{\alpha})^{-1}]^+$.

iv) Predictors arising from (simple) ridge estimators

$$R_{k_t} = \Sigma \ell_i \frac{d_i}{d_i + k_t} \hat{\alpha}_i$$

where $k_t$ is based on the data, i.e., $h_i(\hat{\alpha}) = d_i/(d_i + k_t(\hat{\alpha}))$. $k$'s discussed include:

$k_1(\hat{\alpha}) = \sigma^2 p (\hat{\alpha}^T \hat{\alpha})^{-1}$ (Hoerl, Kennard, and Baldwin 1975),

$k_2(\hat{\alpha}) = \sigma^2 p Z^{-1}$  (Lawless and Wang 1976),

$k_3(\hat{\alpha})$, the solution to $\Sigma \hat{\alpha}_i^2 d_i^2 (d_i + k_3)^{-2} = \Sigma \hat{\alpha}_i^2 - \sigma^2 \Sigma d_i^{-1}$

(McDonald and Galarneau 1975),

$k_4(\hat{\alpha})$, the solution to $\Sigma \hat{\alpha}_i^2 d_i (d_i + k_4)^{-1} = \sigma^2 p$

(the RIDGM estimator of Dempster, Schatzoff and Wermuth 1977).

A subclass of (2) which includes (i), (ii), (iii), and $R_{k_2}$ has the form

$$\delta(\hat{\alpha}) = \Sigma h_i(U,Z) \ell_i \hat{\alpha}_i \ . \tag{3}$$

A further subclass which still includes (i), (ii), and (iii) is

$$\delta(\hat{\alpha}) = h(U,Z) \cdot U. \tag{4}$$

When $D = I$, all of the aforementioned estimators belong to (4).

Taking another point of view (see e.g. Thompson (1968)), if $h_i$ in (3) is constant, the optimal $h_i$ to minimize the MSE are easily obtained:

$$h_i^* = \frac{\theta}{\sigma^2 + \lambda} \frac{\alpha_i}{\ell_i} \ . \tag{5}$$

An estimator of $h_i^*$ would be of the form $c_i(\hat{\alpha}, \sigma^2)$ leading to a predictor belonging to (2). If (5) was estimated by $c(U,Z,\sigma^2) \cdot \hat{\alpha}_i / \ell_i$ the class (4) results.

Suppose we take a Bayesian approach using a prior which centers $\theta$ at 0, where we want to do well. More precisely, let $Q$ be an orthogonal matrix such that $Q D^{\frac{1}{2}} \alpha = \binom{\theta/\sqrt{q}}{\eta}$ where $\eta$ is $(p-1) \times 1$ and $\eta^T \eta = \phi$. If we take as our prior

$$\left(\begin{matrix} \theta/\sqrt{q} \\ \eta \end{matrix}\right) \sim N(0, \left(\begin{matrix} \gamma & 0 \\ 0 & \rho\gamma I_{p-1} \end{matrix}\right)), \ \rho \text{ known},$$

then under squared error loss, the Bayes estimate of $\theta$ is $(\gamma + \sigma^2)^{-1} \cdot \gamma U$. Since $(U,Z)$ is sufficient under the marginal distribution of $\omega = QD^{\frac{1}{2}}\hat{\alpha}$ an "empirical Bayes" estimator of $\theta$ takes the form in (4).

## 3. EXAMINATION OF THE MSE

We can calculate the MSE for the general predictor in (2) in terms of the $h_i$, assuming $\sigma^2$ known.[1]

**Theorem 1.** If $E\left|\frac{\partial h_i}{\partial U} \cdot \hat{\alpha}_i\right| < \infty$, $i = 1,2,\ldots,p$,

$$MSE(\hat{\delta}) = \sigma^2 q + E(\hat{\delta} - U)^2 - 2\sigma^2 E\Sigma \ell_i^2(1 - h_i)$$

$$+ 2\sigma^2 q E\Sigma \ell_i \hat{\alpha}_i \frac{\partial h_i}{\partial U} . \tag{6}$$

**Proof.** By direct calculation

$$MSE(\hat{\delta}) = \sigma^2 q + E(\delta - U)^2 - 2E\{r(\hat{\alpha})(U - \theta)\} \tag{7}$$

where $r(\hat{\alpha}) = \Sigma(1 - h_i)\ell_i \hat{\alpha}_i$. Stein's identity (Stein 1981, p. 1148) converts the right-most term of (7) to $\sigma^2 q E(\frac{\partial r(\hat{\alpha})}{\partial U})$. Simplification yields (6).

$\frac{\partial h_i}{\partial U}$ would be calculated using the transformation $\hat{\alpha} = D^{-\frac{1}{2}}Q^T \omega$ of the previous section. In the case of (3), it can be calculated directly writing $h_i$ as a function of $U$ and $V$. For predictors of the form (4), $MSE(\hat{\delta})$ depends only on $\theta$ and $\phi$ and is given as Corollary 1.

**Corollary 1.** For the predictors in (4), if $E\left|U\frac{\partial h}{\partial U}\right| < \infty$

$$MSE(\delta) = \sigma^2 q + E(1 - h)^2 U^2 + 2\sigma^2 q EU \frac{\partial h}{\partial U} - 2\sigma^2 q E(1 - h). \tag{8}$$

Under (4) choices of $h$ in the literature are such that $h$ is symmetric in $U$ about 0 and restricted to $[0,1]$. Using essentially

4

the argument of Efron and Morris (1976, p. 14) positive part restriction of h uniformly reduces risk. Restriction of $h \le 1$ is less clear. Taking $h > 0$ the predictor $h* \cdot U$ where $h* = \min(h,1)$ does not necessarily dominate $h \cdot U$. For example, let

$$h(U,V) = \begin{cases} 1 + c, & a^2 < U^2 < b^2 \\ 1, & \text{elsewhere} \end{cases}$$ . Then at each $\phi$, for $|\theta|$ suffi-

ciently large, MSE of $h(U,V)U$ is less than MSE of $h*(U,V)$. Nonetheless, to improve in a neighborhood of a specified $\theta_0$ requires convex combinations of $U$ and $\theta_0$. Theorem 2 details MSE properties of predictors in (4) relative to the MSE of $U$.

Theorem 2. For $\delta(\hat{\alpha})$ in (4) with $h \in [0,1]$, let $h$ be symmetric in U about 0. Let $g = (1 - h)U$ with $\lim\sup_{|U| \to \infty} g = 0$ and assume $\frac{\partial g}{\partial U}$

exists for all $U$. Finally, assume that the Lebesgue measure of $A = \{(U,V) : h(U,V) < 1\}$ is greater than 0. Then,

(i) For each $\phi$ there is a neighborhood $N_\phi$ of $\theta = 0$ where $MSE(\delta;\theta,\phi) < \sigma^2 q$.

(ii) $MSE(\delta;\theta,\phi)$ is bounded and $\lim_{|\theta| \to \infty} MSE(\delta;\theta,\phi) = \sigma^2 q$.

(iii) $MSE(\delta;\theta,\phi)$ is symmetric in $\theta$ about 0 and $\frac{\partial MSE(\delta;\theta,\phi)}{\partial \theta}\Big|_{\theta=0} = 0$.

(iv) $g^2 - 2\frac{\partial g}{\partial U}$ changes sign at least once in $0 < U < \infty$. If $g^2 - 2\frac{\partial g}{\partial U}$ changes sign b times in $0 < U < \infty$, then for fixed $\phi$, $MSE(\delta;\theta,\phi) - \sigma^2 q$ changes sign at most 2b times.

Proof. The proof of (i) is clear since $MSE(\delta;0,\phi) < \sigma^2 q$. For (ii),

$$MSE(\delta;\theta,\phi) = \sigma^2 q + Eg^2 - 2E(U - \theta)g. \qquad (9)$$

Given $\epsilon$, $\exists u_0$ such that for all V, $U > u_0 \Rightarrow |g| < \epsilon$ and $\exists \theta_0 > 0$ such that $|\theta| > \theta_0 \Rightarrow P(|U| > u_0) > 1 - \epsilon$. Then the second term and the third term (using the Cauchy-Schwarz Inequality) in (9) can be made arbitrarily small as $|\theta| \to \infty$. It is clear that the r.h.s. of (9) is bounded. (iii) is obvious. The first part of (iv) follows since U is admissible. The second part follows from

5

the sign change theorem of Karlin (1957) by noting that

$$MSE(\hat{\epsilon};\theta,\phi) - \sigma^2 q = E(g^2 - 2\,\frac{\partial g}{\partial U}).$$

**Remark 1.** Predictors in (i), (ii), (iii) of Section 2 satisfy the conditions of Theorem 2.

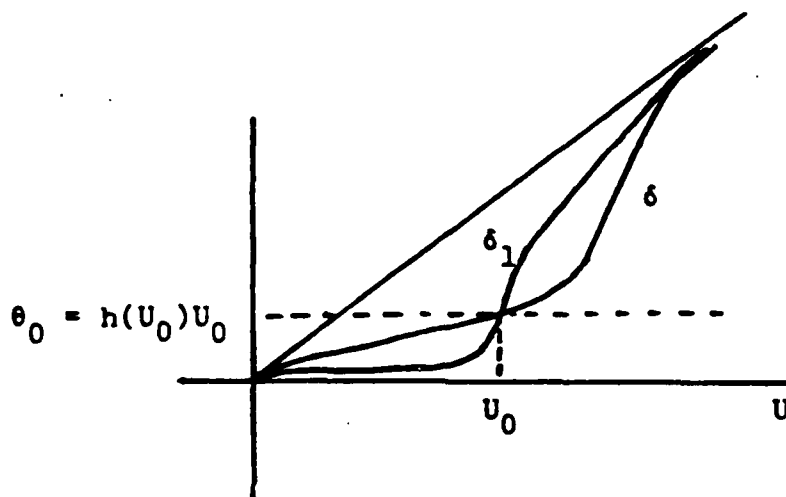**Remark 2.** Result (ii) is a simple case of the "tail mini-maxity" notion of Berger (1976).

**Remark 3.** In (iii), $\inf_{\epsilon} MSE(\delta;\theta,\phi)$ need not occur at $\theta = 0$. If, however, $h(U,V)$ is increasing in $|U|$ it must as may be shown by establishing the result for h, a step function in U. An induction argument proves this.

**Remark 4.** If $b = 1$ in (iv), then a graph of $MSE(\hat{\epsilon};\theta,\phi)$ for $\epsilon \geq 0$ must start below $c^2 q$ at $\theta = 0$, cross above $c^2 q$ at some $\theta$ and then asymptotically return to $\sigma^2 q$ from above. Any predictor satisfying the conditions of Theorem 2 must necessarily perform worse for a set of $\theta$'s near 0 than for a set arbitrarily far away.

**Remark 5.** No immediate extension of Theorem 2 to $\delta(\hat{\alpha})$ as in (3) is available. For an arbitrary member of (3), MSE depends upon $\theta$ and $\eta$ and, even if each $h_i$ meets the "tail minimaxity" condition, need not approach $\sigma^2 q$ as $|\theta| \to \infty$ for fixed $\eta$.

**Remark 6.** Theorem 2 is readily extended to the comparison of any pair of predictors in (4).

We conclude with a comment on admissibility for the above predictors. Within the class of predictors based solely on U, i.e., $h(U)U$, those meeting the conditions of Theorem 2 will either be admissible or if not then improvement cannot be substantial. We employ ideas of Chow and Hwang (1984). Suppose $\delta_1(U)$ is to dominate $\delta_0 = h(U)U$ meeting the conditions of Theorem 2. We can write $\hat{\epsilon}_1$ as $h^*(U)U$, and assume $h^* \geq 0$. For $\delta_1$ to dominate $\hat{\epsilon}_0$ requires, when $|U|$ is large, that generally $h^*$ be closer to 1 than h and that, when $|U|$ is small, generally $h^*$ be closer to 0 than h. A simplified picture of $\delta_0, \delta_1$ for $U > 0$ might look like

$$\theta_0 = h(U_0)U_0$$

But, at $\xi = \xi_0$, it would be almost impossible for $\xi_1$ to dominate. Thus, the simplest $h^*$ which realistically could dominate would have to have at least 3 sign changes for $h - h^*$ on $U > 0$. For such an $h^*$, its form would be complicated, domination would be difficult to show, and improvement would be minimal.

This argument does not extend to the more general class (4). Though $U$ and $V$ are independent, conditioning on $V$ in the above heuristic leads to $\xi_0$ depending upon $V$. We, nonetheless, conjecture "approximate admissibility" for members of (4) meeting the conditions of Theorem 2.

### FOOTNOTE

[1]When $\sigma^2$ is unknown, we customarily assume an estimator $S^2$ of $\sigma^2$ such that $\nu S^2 \sim \sigma^2 \chi_\nu^2$ independent of $\hat{\alpha}$. In the foregoing predictors, $\sigma^2$ is replaced by $cS^2$. As Lawless (1981, pp. 463-464) notes, when $\nu \to \infty$ and even when $\nu$ is moderate, resulting MSE will differ little from that with $\sigma^2$ known.

## BIBLIOGRAPHY

Baranchik, Alvin J. (1964). "Multiple Regression and Estimation
of the Mean of a Multivariate Normal Distribution," _Stanford
University Tech. Report #51_.

Berger, J.O. (1976). "Tail minimaxity in location vector problems
and its applications," _Annals of Statistics_, 4, 33-50.

Chow, M.S. and Hwang, J.T. (1984). "The comparison of estimators
for the noncentrality of a chi-square distribution," Cornell
Statistical Center, _Technical Report_.

Cohen, A. (1965). "Estimates of linear combinations of the para-
meters in the mean vector of a multivariate distribution,"
_Annals of Mathematical Statistics_, 36, 78-87.

Dempster, A.P., Scharzoff, M. and Wermuth, N. (1977). "A simula-
tion study of alternatives to ordinary least squares," _Journal
of the American Statistical Association_, 72, 77-106.

Efron, B. and Morris, C. (1976). "Families of minimax estimators of
the mean of a multivariate normal distribution," _Annals of
Statistics_, 4, No. 1, 11-21.

Hoerl, A.E., Kennard, R.W. and Baldwin, K.F. (1975). "Ridge
regression: Some simulations," _Communications in Statistics_, 4,
105-123.

Karlin, S. (1957). "Polya type distributions II," _Annals of
Mathematical Statistics_, 28, No. 2, 281-308.

Lawless, J.F. (1981). "Mean squared error properties of general-
ized ridge estimators," _Journal of the American Statistical
Association_, 76, 462-466.

Lawless, J.F. and Wang, P. (1976). "A simulation study of ridge
and other regression estimators," _Communications in Statistics_,
A, 5, 307-323.

McDonald, G.C. and Galarneau, D.L. (1975). "A Monte Carlo evalu-
ation of some ridge type estimators," _Journal of the American
Statistical Association_, 70, 407-416.

Mehta, J.S. and Srivastava, R. (1971). "Estimation of the mean by
shrinking to a point," _Journal of the American Statistical
Association_, 66, 86-91.

Radhakrishnan, R. (1970). "Some Estimation Problems in Multivariate Analysis," unpublished Ph.D. Thesis, Carnegie Mellon University.

Sclove, S.L. (1968). "Improved estimators for coefficients in linear regression," Journal of the American Statistical Association, 63, 596-606.

Stein, C. (1981). "Estimation of the mean of a multivariate normal distribution," Annals of Statistics, 9, No. 6, 1135-1151.

Thompson, J.R. (1968). "Some shrinkage techniques for estimating the mean," Journal of the American Statistical Association, 63, 113-122.

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|---|---|---|
| 1. REPORT NUMBER<br><br>427 | 2. GOVT ACCESSION NO. | 3. RECIPIENT'S CATALOG NUMBER |
| 4. TITLE (and Subtitle)<br><br>Mean Square Error Behavior For Prediction In Linear Regression Models | | 5. TYPE OF REPORT & PERIOD COVERED<br><br>TECHNICAL REPORT |
| | | 6. PERFORMING ORG. REPORT NUMBER |
| 7. AUTHOR(s)<br><br>Alan E. Gelfand | | 8. CONTRACT OR GRANT NUMBER(s)<br><br>N00014-89-J-1627 |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS<br><br>Department of Statistics<br>Stanford University<br>Stanford, CA 94305 | | 10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS<br><br>NR-042-267 |
| 11. CONTROLLING OFFICE NAME AND ADDRESS<br><br>Office of Naval Research<br>Statistics & Probability Program Code 1111 | | 12. REPORT DATE<br>March 7, 1990 |
| | | 13. NUMBER OF PAGES<br>11 |
| 14. MONITORING AGENCY NAME & ADDRESS(If different from Controlling Office) | | 15. SECURITY CLASS. (of this report)<br><br>UNCLASSIFIED |
| | | 15a. DECLASSIFICATION/DOWNGRADING SCHEDULE |

16. DISTRIBUTION STATEMENT (of this Report)

APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNLIMITED

17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, If different from Report)

18. SUPPLEMENTARY NOTES

19. KEY WORDS (Continue on reverse side If necessary and Identify by block number)

prediction; adaptive estimators; linear regression; mean square error.

20. ABSTRACT (Continue on reverse side If necessary and Identify by block number)

For the problem of individual prediction in linear regression models, that is, estimation of a linear combination of regression coefficients, mean square error behavior of a general class of adaptive predictors is examined.